# Diffusion-Based Scene Graph to Image Generation with Masked Contrastive Pre-Training

Ling Yang[1*]    Zhilin Huang[2*]    Yang Song[3]    Shenda Hong[1]    Guohao Li[4]    Wentao Zhang[5]
Bin Cui[1]    Bernard Ghanem[4]    Ming-Hsuan Yang[6,7]
[1]Peking University    [2]Tsinghua University    [3]OpenAI    [4]KAUST
[5]Mila    [6]University of California, Merced    [7]Google Research
yangling0818@163.com, huang-zl22@mails.tsinghua.edu.cn, songyang@openai.com,
wentao.zhang@mila.quebec, {hongshenda, bin.cui}@pku.edu.cn,
{guohao.li, bernard.ghanem}@kaust.edu.sa, mhyang@ucmerced.edu

## Abstract

*Generating images from graph-structured inputs, such as scene graphs, is uniquely challenging due to the difficulty of aligning nodes and connections in graphs with objects and their relations in images. Most existing methods address this challenge by using scene layouts, which are image-like representations of scene graphs designed to capture the coarse structures of scene images. Because scene layouts are manually crafted, the alignment with images may not be fully optimized, causing suboptimal compliance between the generated images and the original scene graphs. To tackle this issue, we propose to learn scene graph embeddings by directly optimizing their alignment with images. Specifically, we pre-train an encoder to extract both global and local information from scene graphs that are predictive of the corresponding images, relying on two loss functions: masked autoencoding loss and contrastive loss. The former trains embeddings by reconstructing randomly masked image regions, while the latter trains embeddings to discriminate between compliant and non-compliant images according to the scene graph. Given these embeddings, we build a latent diffusion model to generate images from scene graphs. The resulting method, called SGDiff, allows for the semantic manipulation of generated images by modifying scene graph nodes and connections. On the Visual Genome and COCO-Stuff datasets, we demonstrate that SGDiff outperforms state-of-the-art methods, as measured by both the Inception Score and Fréchet Inception Distance (FID) metrics. We will release our source code and trained models at https://github.com/YangLing0818/SGDiff.*
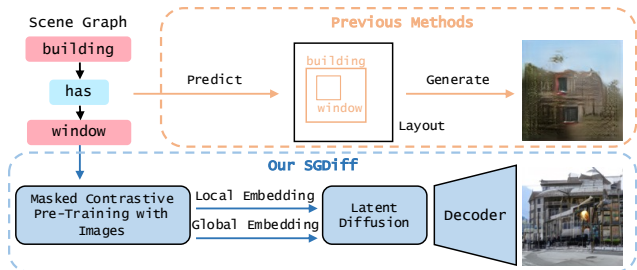
---

*Contributed equally.

Figure 1. **SGDiff *vs*. Previous methods.** Instead of relying on manually specified scene graph representations such as scene layouts, we pre-train our scene graph embeddings with masked autoencoding loss and contrastive learning, explicitly maximizing their predictive power for graph-image alignment. Conditioned on such embeddings, our latent diffusion model *SGDiff* outperforms prior works on scene graph to image generation.

## 1. Introduction

Image generation has made remarkable progress in the past few years [4, 9, 12, 40], largely due to the success of diffusion and score-based generative models [19, 48–50]. These methods allow for the creation of realistic and diverse image samples [11, 39, 45], which users can specify through various forms—labels [11, 20, 25, 50], captions [42, 45], segmentation masks [8], sketches [60], stroke paintings [33], and more [62]. However, these types of specifications often fall short when it comes to complex relations between multiple objects in images. Instead, scene graphs provide a concise and accurate way of depicting objects and their relations to one another [23, 29]. It is therefore crucial to investigate image generation based on scene graphs as a means of synthesizing complex scenes [22].

A key challenge in generating images from scene graphs is ensuring that the resulting image closely aligns with the

input scene graph. To this end, generative models must be able to understand the correspondence between the two vastly different data domains: images and graphs. Existing methods [1, 17, 22, 30] mainly address this challenge by using an image-like representation of scene graphs, often in the form of scene layouts, to create coarse sketches for guiding the image generation process. These sketches are then refined by generative models to produce realistic images that follow the specifications given by the scene graph.

While intermediate representations such as scene layouts can be useful, they are often crafted manually and are not specifically designed to facilitate the alignment between images and graphs. For instance, in the case of scene layouts, nodes in scene graphs are usually mapped to bounding boxes and connections are mapped to their spatial layouts. However, not all connections within scene graphs can be accurately translated to spatial layouts, such as `eating` and `looking at`. Additionally, some relations, such as `behind`, `inside`, and `in front of`, all correspond to similar spatial relations in scene layouts, creating ambiguity. These intermediate scene layout representations may also contain extraneous information that complicates the training of downstream generative models.

To overcome these limitations, we propose learning intermediate representations that explicitly maximize the alignment between scene graphs and images. We provide an illustration of our approach in Fig. 1. Specifically, we pre-train a scene graph encoder on graph-image pair datasets to produce embeddings that extract both local and global information from scene graphs, while maximizing their alignment with images. To extract local information, we introduce a masked autoencoding loss, randomly masking out objects in an image and reconstructing the missing portion using unmasked regions and embeddings acquired from the encoder. To gain global information, we leverage contrastive learning to train our encoder to discern between images that do and do not adhere to scene graphs. By combining embeddings obtained from both approaches, we obtain compact intermediate representations of scene graphs that facilitate the alignment between graphs and images.

We showcase the effectiveness of our scene graph embeddings by building a latent diffusion model [42, 56] that generates images from scene graphs with the aid of our pre-trained embeddings. We evaluate the importance of local and global embeddings through ablation studies, and demonstrate the clear advantages our approach has over traditional intermediate layout representations. Our model, dubbed *SGDiff*, successfully generates images that capture accurate local and global structures of scene graphs. Additionally, our model enables the semantic manipulation of images through scene graph surgery. We evaluate SGDiff on standard datasets such as Visual Genome (VG) [29] and COCO-Stuff [5], and find that it performs better than cur-rent state-of-the-art approaches in both qualitative comparison and quantitative measurements.

## 2. Related Work

**Diffusion Models for Conditional Image Synthesis.** Diffusion models [62] generate data samples by learning to reverse a prescribed diffusion process that converts data to noise. First introduced by Sohl-Dickstein *et al*. [48] and later improved by Song & Ermon [49] and Ho *et al*. [19], they are now able to generate image samples with unprecedented quality and diversity [13, 44, 45]. Diffusion models excel at conditional image synthesis using various forms of user guidance, such as text and images. Existing conditional diffusion models often leverage auxiliary classifiers [11, 20, 25, 34, 50] to incorporate conditional information into the data generating process, using methods of classifier guidance [11, 50] or classifier-free guidance [20]. Latent Diffusion Models (LDMs) [42, 56] reduce the training cost for high resolution images by learning the diffusion model in a low-dimensional latent space. They also incorporate conditional information into the sampling process via cross attention [58]. Similar techniques are employed in DALLE-2 [39] for image generation from text, where the diffusion model is conditioned on text embeddings obtained from CLIP latent codes [37]. Alternatively, Imagen [45] implements text-to-image generation by conditioning on text embeddings acquired from large language models (*e.g*., T5 [38]). Despite all this progress on diffusion-based conditional image synthesis, generating images from graph-structured data is under-explored. We fill this gap by designing the first diffusion model for image generation from scene graphs, leveraging new scene graph embeddings constructed from self-supervised learning.

**Image Generation from Scene Graphs.** Scene graphs are graph-structured data for describing multiple objects and their complex relationships in scene images, wherein nodes represent objects and edges represent relations [23, 29]. Image generation from scene graphs [22, 55] requires the generative model to reason over both objects and their relations. The first generative model of this kind, Sg2Im [22], proposes a two-stage generation pipeline. First, an embedding model is trained to map scene graphs to scene layouts, which are image-like representations that capture the coarse structure of images to generate. Second, a generative model is trained to refine scene layouts into realistic images. Most subsequent work on this task follows the same pipeline. For example, WSGC [17] accounts for semantic equivalence in graph representations by canonicalizing scene graphs before mapping them to scene layouts. Li *et al*. [30] and Ashual & Wolf [1] leverage a repository of external reference images to improve the quality of scene layouts. Other works that rely on scene layout style representations include [16, 31, 52, 53, 64]. In lieu of manually crafted
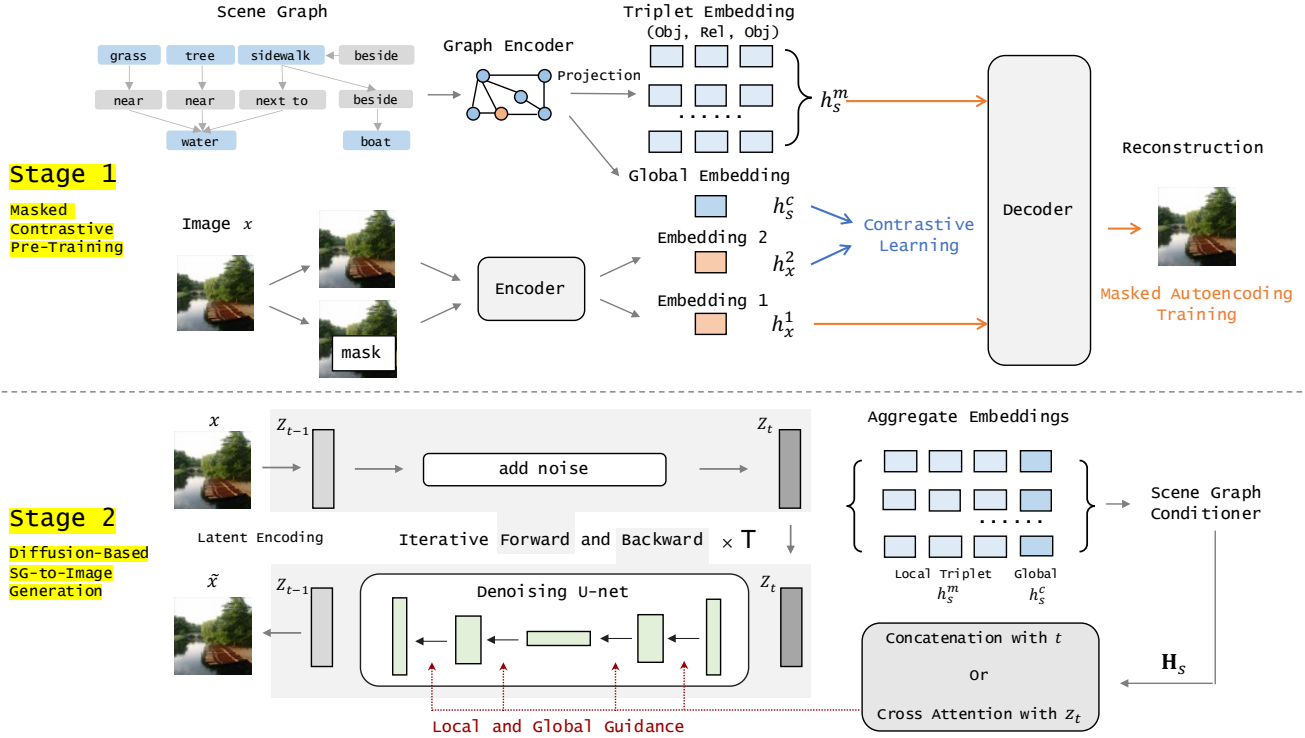
Figure 2. **Schematic illustration of SGDiff.** In the first stage, we pre-train a scene graph encoder using the combination of a masked autoencoding loss and a contrastive loss. In the second stage, we build a latent diffusion model that conditions on embeddings produced by the scene graph encoder.

scene layouts, we propose learning scene graph embeddings that are both concise and predictive of graph-image alignment. We then use these embeddings to build a latent diffusion model for scene graph to image generation, avoiding the limitations of scene layouts.

## 3. Method

In what follows, we first discuss how to learn effective embeddings of scene graphs via self-supervised learning, then focus on using these embeddings to build diffusion models for scalable scene graph to image generation.

### 3.1. Masked Contrastive Pre-Training

Scene layouts are manually constructed representations of scene graphs (SGs). While effective in many cases, they are not specifically optimized to capture all the necessary information from SGs for generating the corresponding scene images. This can lead to suboptimal alignment between SG inputs and generated images. To overcome this limitation, we propose to directly learn such SG representations via self-supervised learning. Specifically, given a dataset that contains SG-image pairs, we learn an SG encoder to pro-

duce embeddings that capture both local and global information of the input SG, while maximizing their alignment with the corresponding scene images.

**Notations.** Given a set of objects $\mathcal{C}_o$ and a set of relations $\mathcal{C}_r$, we denote a scene graph $s$ with a tuple $(O, \mathcal{R})$. Here $O = \{o_i \in \mathcal{C}_o\}_{i=1}^n$ represents the set of objects in the scene and $\mathcal{R} = \{r_{ij} \in \mathcal{C}_r\}_{1 \leq i,j \leq n}$ denotes the set of relations that exist between these objects. We use a triplet $(o_i, r_{ij}, o_j)$ to denote a directed connection from $o_i$ to $o_j$, representing the relation tuple (subject, predicate, object). We denote by $\mathcal{N}_{out}(o_i)$ (resp. $\mathcal{N}_{in}(o_i)$) the set of children (resp. parents) for node $o_i$. To generate an effective representation for $s$, we embed both objects and relations by iterating the following:

$$h_{o_i} = \text{Pool}(\{f_o^{out}(h_{o_i}, h_{r_{ij}}, h_{o_j})\}_{j \in \mathcal{N}_{out}(o_i)}$$
$$\cup f_o^{in}(h_{o_j}, h_{r_{ji}}, h_{o_i})\}_{j \in \mathcal{N}_{in}(o_i)}), \quad (1)$$
$$h_{r_{ij}} = f_r(h_{o_j}, h_{r_{ij}}, h_{o_i}), \quad (2)$$

where $h_{o_i} \in \mathbb{R}^{d_o}$ is the embedding of object $o_j$, $h_{r_{ij}} \in \mathbb{R}^{d_r}$ is the embedding of connection $r_{ij}$ (the relation from $o_i$ to $o_j$), and $f_o^{out}$, $f_o^{in}$, $f_r$ denote separate graph convolutional layers [28, 47]. Here $\text{Pool}(\{\})$ denotes the average pooling

operator. To train these object and relation embeddings, we oftentimes need an image encoder. The embedding from this encoder is denoted as $h_x$ for image $x$.

Below, we introduce two self-supervised techniques for learning object and relation embeddings, focusing on extracting local and global information from SG inputs.

**Pretraining with Masked Autoencoding.** Masked pretraining is a preeminent technique in image/text representation learning [3, 6, 10, 14, 61] and visual-language modeling [32, 51, 63]. Inspired by its success in these applications, we propose to use a similar technique for learning our SG embeddings.

In particular, we randomly choose a triplet $(o_i, r_{ij}, o_j)$ in the scene graph $s$ and mask out objects $o_i$ and $o_j$ in the corresponding scene image $x$. We denote this masked area as $x_m$, and the remaining image as $x_{\setminus m}$. To train the SG encoder, we consider the task of predicting $x_m$ from $x_{\setminus m}$ and embeddings obtained from the SG encoder. Specifically, we concatenate object and relation embeddings from the SG encoder to form $h_s^m$, defined as:

$$h_s^m = \texttt{concat}(\{(f_o^m(h_{o_i}), f_r^m(h_{r_{ij}}), f_o^m(h_{o_j}))\}), \quad (3)$$

where $f_o^m : \mathbb{R}^{d_o} \rightarrow \mathbb{R}^d$, $f_r^m : \mathbb{R}^{d_r} \rightarrow \mathbb{R}^d$ are MLPs with one hidden layer and ReLU activation functions, and $\texttt{concat}$ stacks all triplets of the form $(f_o^m(h_{o_i}), f_r^m(h_{r_{ij}}), f_o^m(h_{o_j}))$ in the SG to create the vector $h_s^m$. We jointly train the SG encoder and an auxiliary decoding model $d_\theta$ to minimize the following masked autoencoding loss:

$$\mathcal{L}_{\text{masked}} = \mathbb{E}_{(s,x)\sim D} \|x_m - d_\theta(x_{\setminus m}, h_s^m)\|_2^2, \quad (4)$$

where $(s, x)$ is sampled uniformly at random from $D$, a dataset of graph-image pairs. By training the SG embeddings to reconstruct randomly masked areas in scene images, we explicitly encourage the graph embeddings to encode local structural information that focuses on predicting fine-grained image details from scene graphs.

**Pretraining with Contrastive Learning.** Contrastive pre-training [7, 15] is a widely adopted technique for learning shared representations across multiple data modalities. They have found success in many visual-language modeling applications [36, 37]. We propose to leverage contrastive learning as a second way to train our SG embeddings, with a focus on capturing global structural information. Specifically, we compute a graph-level embedding $h_s^c$ by concatenating all object and relation embeddings obtained from the SG encoder, that is,

$$h_s^c = f^c(\texttt{concat}(\text{Pool}(\{h_o\}_{o\in\mathcal{O}}), \text{Pool}(\{h_r\}_{r\in\mathcal{R}}))), \quad (5)$$

where $f^c : \mathbb{R}^{d_r+d_o} \rightarrow \mathbb{R}^d$ is an MLP with one hidden layer and ReLU activation functions. We call $(s, x^+)$ a positive pair if $x^+$ complies with $s$, and $(s, x^-)$ a negative pair if $x^-$ does not adhere to $s$. Positive pairs can be sampled from the graph-image pair dataset $D$, whereas negative pairs are generated by uniformly choosing an image $x^-$ in $D$ that does not match $s$. We optimize the following cross entropy [35] loss to learn graph embeddings that can discern positive pairs from negative ones:

$$\mathcal{L}_{\text{contrastive}}(f; \tau, k) =$$
$$\mathbb{E}_{\substack{(s,x^+)\sim D \\ \{x_i^-\}_{i=1}^k \sim D_s}} \left[ -\log \frac{\exp(h_s^{c\top} h_{x^+}/\tau)}{\exp(h_s^{c\top} h_{x^+}/\tau) + \sum_i \exp(h_s^{c\top} h_{x_i^-}/\tau)} \right], \quad (6)$$

where $\tau$ is a learnable multiplicative scalar that acts as a temperature parameter, $h_x$ denotes the embeddings of an image $x$ produced by a trainable image encoder, and $D_s$ denotes the set of images in dataset $D$ that do not comply with the SG $s$. This contrastive training objective optimizes our SG embeddings to capture global structures that can identify whether images are in line with the SGs or not.

We combine both the masked autoencoding loss and the contrastive loss to form our training objective for SG embeddings:

$$\mathcal{L} = \mathcal{L}_{\text{masked}} + \lambda \mathcal{L}_{\text{contrastive}}, \quad (7)$$

where $\lambda > 0$ is a hyperparameter. With this objective function, we can train an effective SG encoder that converts SGs to embeddings without losing predictive information of the matching scene images. Such embeddings provide strong conditioning signals that facilitate downstream diffusion models to generate images from SGs.

### 3.2. Diffusion-Based SG to Image Generation

In diffusion-based conditional image synthesis, we train diffusion models to sample from $p(x \mid y)$, where $x$ is an image, and $y$ is a conditioning signal, typically taking the form of texts in text-to-image generation [42, 45] and image editing [2, 24, 57], or reference images in the context of image translation [33]. We consider the task of image generation from SGs in this work. With the SG encoder obtained from the previous section, we can build diffusion models to generate images from SGs by setting $y$ to the SG embeddings. For scalable image modeling, we build upon the framework of latent diffusion [42, 56], where the diffusion model is trained in a low-dimensional latent space obtained from a pre-trained autoencoder.

**Diffusion in the Latent Space.** We train our diffusion models in a low-dimensional latent space in order to improve the computational efficiency for generating high-

resolution images. At the first step, we pre-train an autoencoder to map high-dimensional images to low dimensional latent codes. Specifically, the encoder $f_{enc}$ is trained to transform the image $x \in \mathbb{R}^{H \times W \times 3}$ into a latent code $z \in \mathbb{R}^{h \times w \times c}$ with a downsampling factor $k = H/h = W/w$, and the decoder is trained to reconstruct $x$ from $z$. To regularize the distribution of $z$ for stable generative modeling, we additionally minimize the KL divergence from the distribution of $z$ to a standard Gaussian distribution, in the same vein as Variational Autoencoders [27, 41].

After training both the encoder and the decoder, we use the encoder to generate latent codes for all images in the training dataset, then train a diffusion model on these latent codes separately. In particular, given the latent code $z$ for a randomly sampled training image $x$, we convert it to noise with a Markov process defined by the transition kernel $q(z_t \mid z_{t-1}) = \mathcal{N}(z_t; \sqrt{\alpha_t}z_{t-1}, (1 - \alpha_t)\mathbf{I})$, where $t = 1, 2, \cdots, T$, $z_0 = z$, and $\alpha_t$ is a hyper-parameter that controls the rate of noise injection. When the amount of noise is sufficiently large, $z_T$ becomes approximately distributed according to $\mathcal{N}(0, \mathbf{I})$. In order to convert noise back to data for sample generation, we have to estimate the reverse diffusion process by learning the reverse transition kernel $p_\theta(z_{t-1} \mid z_t) = \mathcal{N}(\mu_\theta(z_t), \Sigma_\theta(z_t))$ as an approximation to $q(z_{t-1} \mid z_t)$. Following Ho *et al.* [19], we parameterize $\mu_\theta(z_t)$ with a neural network $\epsilon_\theta(z_t, t)$ (called the score model [49, 50]) and fix $\Sigma_\theta$ to be a constant. The score model can be optimized with denoising score matching [21, 59]. For sample generation, we first generate a latent code $z$ with the diffusion model, then produce an image sample $x$ through the pre-trained decoder.

**Conditioning on SG Embeddings.** With the method in Sec. 3.1, we can learn an SG encoder to produce two types of scene graph embeddings, $h_s^m$ and $h_s^c$, defined according to Eq. (3) and Eq. (5). To combine these embeddings, we first merge $h_s^c$ and $h_s^m$ by summing $h_s^c$ with each triplet of the form $\left(f_o^m(h_{o_i}), f_r^m(h_{r_{ij}}), f_o^m(h_{o_j})\right)$ (*cf*., Eq. (3)), which has been generated in the process of computing $h_s^m$. This gives us a new embedding for each relation:

$$h_{r_{ij}}^{sum} = f_o^m(h_{o_i}) + f_r^m(h_{r_{ij}}) + f_o^m(h_{o_j}) + h_s^c. \quad (8)$$

Afterwards, we concatenate and transform them to get our final embedding, given by

$$\mathbf{H}_s = \psi_{cond}(\texttt{concat}(\{h_{r_{ij}}^{sum}\}_{r_{ij} \in \mathcal{R}})), \quad (9)$$

where $\psi_{cond}$ is a trainable model called the SG conditioner. We then use the resulting embedding $\mathbf{H}_s$ to guide the generation of our latent diffusion model.

Specifically, we incorporate the embedding into each block of the UNet [43] backbone of the score model in latent diffusion. We experiment with two different conditioning methods. The first is concatenating the SG embedding

$\mathbf{H}_s \in \mathbb{R}^{N \times d_s}$ with time embedding $t$, in which case the conditional score model $\epsilon_\theta$ takes the form of:

$$\epsilon_\theta(z_t, t; \mathbf{H}_s) = \epsilon_\theta(z_t, \texttt{concat}(t, \mathbf{H}_s)). \quad (10)$$

The second is applying cross-attention [58] to combine the SG embedding $\mathbf{H}_s \in \mathbb{R}^{N \times d_s}$ with the noisy latent code $z_t$, where $\epsilon_\theta$ is given by

$$\epsilon_\theta(z_t, t; \mathbf{H}_s) = \epsilon_\theta(\texttt{Cross-Attention}(z_t, \mathbf{H}_s), t). \quad (11)$$

Here $z_t' = \texttt{Cross-Attention}(z_t, \mathbf{H}_s)$ is defined as:

$$z_t' = \texttt{softmax}\left(\frac{(W_Q \cdot \phi(z_t))(W_K \cdot \mathbf{H}_s)^\top}{\sqrt{d}}\right) \cdot (W_V \cdot \mathbf{H}_s), \quad (12)$$

where $W_Q \in \mathbb{R}^{d \times d_{z_t}}$, $W_K \in \mathbb{R}^{d \times d_{\mathbf{H}_s}}$ and $W_V \in \mathbb{R}^{d \times d_{\mathbf{H}_s}}$ are learnable matrices, and $\phi(\cdot)$ is a learnable neural network. We can directly train these conditional score models to obtain our latent diffusion model, dubbed *SGDiff*.

## 4. Experimental Results

**Datasets and Evaluation Metrics.** The Visual Genome (VG) [29] dataset contains 108,077 scene graph & image pairs, with additional annotations such as bounding boxes and object attributes. Scene graphs in this dataset have 179 object identities, 80 attributes, and 49 relations. For a fair comparison, we follow previous works [17, 22] to use the standard dataset splits: 62,565 pairs for training, 5,506 for validation, and 5,088 for test. COCO-Stuff [5] contains pixel-wise annotations with 40,000 training images and 5,000 validation images with corresponding bounding boxes and segmentation masks. It has 80 item categories and 91 stuff categories. Following [22], we use synthesized scene graphs and standard dataset splits: 25,000 for training, 1,024 for validation, and 2,048 for test.

We report results with two widely used evaluation metrics. One is Inception Score (IS) [46], which measures both the quality and diversity of synthesized images. Same as previous work, we employ a pre-trained inception network [54] to obtain network activations for computing IS. The IS is better when larger. The other evaluation metric is Fréchet Inception Distance (FID) [18], which is reported to align well with human evaluation. It measures the distance between the distribution of the generated images and that of the real test images, both modeled as multivariate Gaussians. The FID values are better when lower.

**Baselines and Implementation Details.** In our experiments, we choose four previous methods on scene graph to image generation as our baselines: Sg2Im [22], WSGC [17], SOAP [1], and PasteGAN [30]. We follow their evaluation settings in all experiments. For masked contrastive pretraining, we train models using the Adam optimizer [26]
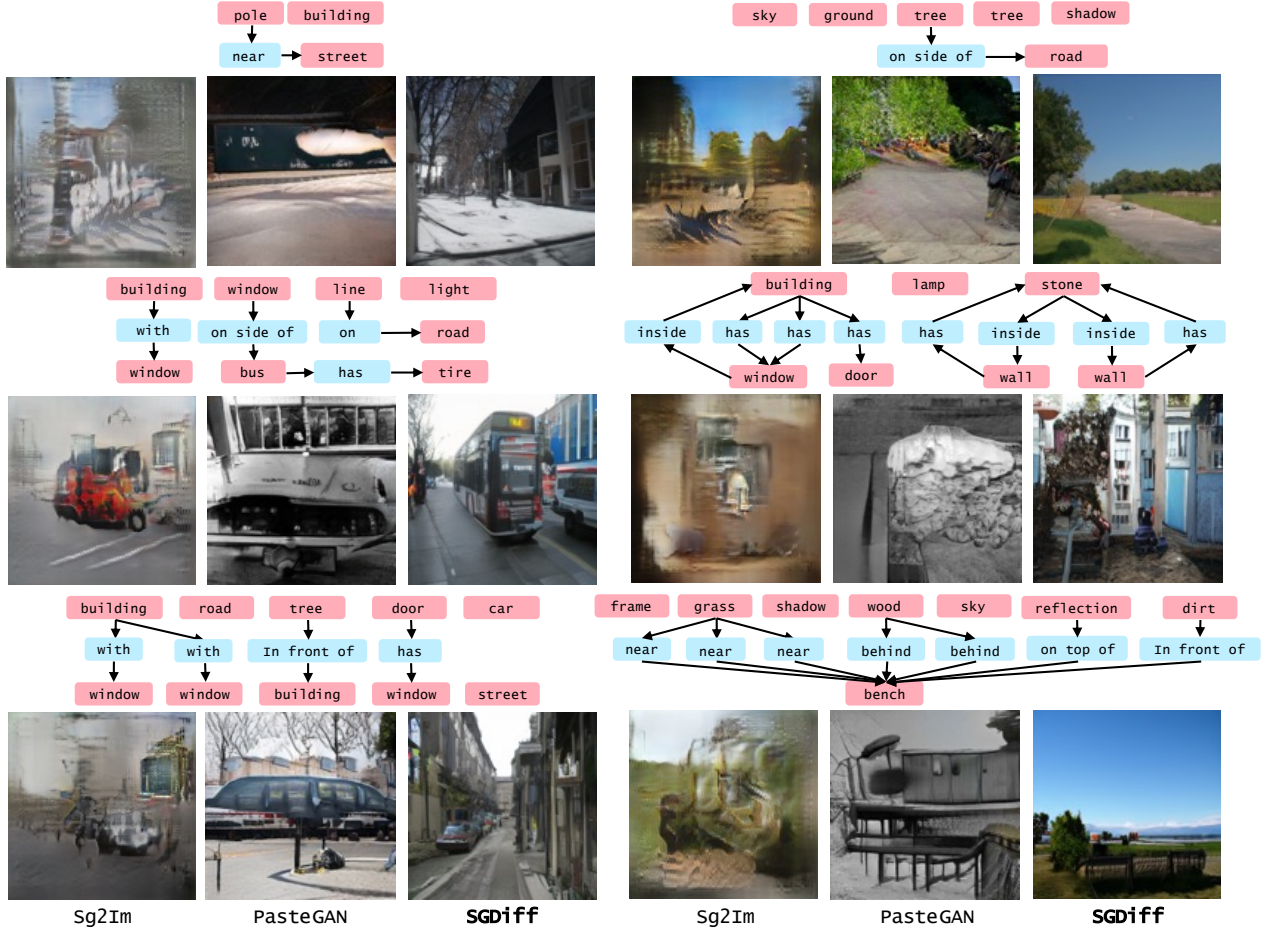
Figure 3. Image samples (128×128) generated by Sg2Im [22], PasteGAN [30], and our SGDiff given the same scene graphs.

with a learning rate of 5e-4 and a batch size of 64 for 100,000 iterations. For latent diffusion training, we train models from scratch using the same optimizer but with a learning rate of 1e-6 and a batch size of 16 for 700,000 iterations. Please consult the Appendix A for more details.

## 4.1. Quantitative Comparisons

Similar to prior works [1, 17, 22, 30], we perform quantitative comparisons on the COCO-Stuff and VG datasets for image resolution 64×64, 128×128, and 256×256. All IS and FID results are provided in Tab. 1. We observe that our SGDiff consistently outperforms existing methods on both evaluation metrics by a significant margin, demonstrating superiority in both generation fidelity and diversity. One possible explanation is that conventional intermediate representations of scene graphs used in previous methods, as exemplified by scene layouts, do not specifically optimize the semantic alignment between scene graphs and images. By contrast, SGDiff directly learns scene graph embeddings, maximizing both local and global semantic compliance with images in a self-supervised manner. Since our

latent diffusion model incorporates these optimized embeddings to guide the generation process, we naturally achieve improved quality in image generation. We verify this hypothesis with rigorous ablation studies in Sec. 4.3. We also test the two methods for conditioning on scene graph embeddings, and find that concatenation with the noisy latent code $z_t$ consistently outperforms concatenation with time embedding $t$.

## 4.2. Qualitative Evaluations

To explore SGDiff's ability of generating images that contain multiple objects and complex relations, we visualize some typical examples (with resolution 128×128) in Fig. 3, where images are placed in ascending order of scene graph complexity. We compare SGDiff with a classical algorithm Sg2Im [22] and the state-of-the-art method PasteGAN [30]. It is clear that images synthesized by SGDiff are not only more realistic, but also comply with the corresponding scene graphs better than both Sg2Im and PasteGAN. In contrast, images generated by Sg2Im are fuzzy and lack fine-grained details. Although PasteGAN can generate

| Method | Inception Score ↑ | | FID ↓ | |
| --- | --- | --- | --- | --- |
| | COCO | VG | COCO | VG |
| Real Img (64×64) | 16.3 ± 0.4 | 13.9 ± 0.5 | - | - |
| Sg2Im [22] | 6.7 ± 0.1 | 5.5 ± 0.1 | 82.8 | 71.3 |
| WSGC [17] | 5.6 ± 0.1 | 8.0 ± 1.1 | 91.3 | 45.3 |
| SOAP [1] | 7.9 ± 0.2 | - | 65.3 | - |
| PasteGAN [30] | 9.1 ± 0.2 | 6.9 ± 0.2 | 50.9 | 58.5 |
| **SGDiff** (with $t$) | **10.6 ± 0.4** | **8.9 ± 0.5** | **26.8** | **27.5** |
| **SGDiff** (with $z_t$) | **11.4 ± 0.4** | **9.3 ± 0.2** | **22.4** | **16.6** |
| Real Img (128×128) | 24.2 ± 0.9 | 17.4 ± 1.1 | - | - |
| Sg2Im [22] | 7.1 ± 0.2 | 6.1 ± 0.1 | 93.3 | 82.7 |
| WSGC [17] | 5.1 ± 0.3 | 7.2 ± 0.3 | 108.6 | 80.4 |
| SOAP [1] | 10.4 ± 0.4 | - | 75.4 | - |
| PasteGAN [30] | 11.1 ± 0.7 | 7.6 ± 0.7 | 70.7 | 61.2 |
| **SGDiff** (with $t$) | **13.1 ± 0.4** | **9.5 ± 0.5** | **32.7** | **29.6** |
| **SGDiff** (with $z_t$) | **14.6 ± 0.9** | **11.4 ± 0.5** | **30.2** | **20.1** |
| Real Img (256×256) | 30.7 ± 1.2 | 27.3 ± 1.6 | - | - |
| Sg2Im [22] | 8.2 ± 0.2 | 7.9 ± 0.1 | 99.1 | 90.5 |
| WSGC [17] | 6.5 ± 0.3 | 9.8 ± 0.4 | 121.7 | 84.1 |
| SOAP [1] | 14.5 ± 0.7 | - | 81.0 | - |
| PasteGAN [30] | 12.3 ± 1.0 | 8.1 ± 0.9 | 79.1 | 66.5 |
| **SGDiff** (with $t$) | **16.0 ± 0.9** | **13.6 ± 0.7** | **40.1** | **36.4** |
| **SGDiff** (with $z_t$) | **17.8 ± 0.8** | **16.4 ± 0.3** | **36.2** | **26.0** |

Table 1. **Inception Scores and FIDs on COCO-Stuff and VG datasets.** Results of previous methods are either directly taken from their original papers or obtained by running official open-source implementations.

images with more details compared to Sg2Im, it tends to miss some important relations specified by the scene graphs, leading to obfuscated results. Both methods are prone to generating images with wrong objects or relations. In contrast, our SGDiff consistently generates realistic images that match scene graphs well, since we rely on scene graph embeddings that are explicitly optimized for local and global semantic alignment. The generated objects have clear object boundaries and more visual details. We place more image samples in Appendix B.

**Semantic Image Manipulation.** To demonstrate the semantic consistency between generated images and scene graphs, we apply our SGDiff to manipulate image samples by modifying objects and relations in scene graph inputs. As shown in Fig. 4, SGDiff can not only produce compliant manipulation results (of resolution 256×256) with respect to objects and relations, but also synthesize perceptually diverse images when conditioned on the same scene graph. The results demonstrate that our model can effectively leverage the scene graph embeddings learned through masked contrastive pre-training.

## 4.3. Ablation Studies

Key to our approach is learning scene graph embeddings via masked contrastive pre-training. Here we perform ablation studies to understand the importance of masked autoencoding loss and contrastive loss in training our scene graph encoder. We also empirically verify that our scene graph embeddings outperform manually crafted scene layout representations when combined with the same latent diffusion model on scene graph to image generation.

**Retrieval Tasks.** We first evaluate the impacts of masked autoencoding loss and contrastive loss on cross-modal semantic alignment through graph-to-image and image-to-graph retrieval experiments. In graph-to-image retrieval, we use scene graph embeddings and image embeddings to search the most semantically similar image for a given scene graph, then report the accuracy of finding the correctly paired image from the given dataset. The image-to-graph retrieval task is defined analogously. All results are provided in Tab. 2. Here "obj." stands for experiments where we only train object embeddings, whereas "Obj. + Rel." represents settings where we train both object and relation embeddings. We observe that using both embeddings boost the performance on all retrieval tasks. With only the contrastive loss, we can already obtain over 70% accuracy in graph-to-image and image-to-graph retrieval. Adding masked auto-encoding loss further improves the accuracy, demonstrating better graph-image alignment.

| Average Accuracy | Graph-to-Image | Image-to-Graph |
| --- | --- | --- |
| Obj. | 68.6 | 69.8 |
| Obj. + Rel. | 70.3 (**+1.7**) | 70.6 (**+0.8**) |
| Contrastive | 70.3 | 70.6 |
| Contrastive + Masked | 73.4 (**+3.1**) | 74.1 (**+3.5**) |

Table 2. Ablation study of masked contrastive pre-training on retrieval tasks.

**Generation from Scene Graphs.** To understand the role of masked contrastive pre-training in scene graph to image generation, we train the same latent diffusion model with different scene graph embeddings: naïve one-hot embeddings, layout embeddings, masked embeddings, contrastive embeddings, and combination of the last two. As for one-hot embeddings, we fix object and relation embeddings as one-hot vectors of their categories, then concatenate these embeddings together to condition the latent diffusion model. For conventional scene layout embeddings, we follow the same settings in PasteGAN [30]. In Tab. 3, we report IS and FID scores for image samples of resolution 256×256. We observe that embeddings obtained from

Figure 4. Semantic image manipulation (256×256) with SGDiff. Red and blue boxes denote object and relation modifications respectively.

either masked pretraining or contrastive pretraining outperforms one-hot embeddings or scene layout representations by a significant margin, and combining them together further improves sample quality.

| Embeddings | Inception Score ↑ | FID ↓ |
|---|---|---|
| One-Hot | $10.1 \pm 1.3$ | 87.1 |
| Layout | $12.3 \pm 1.0$ | 52.7 |
| Masked (Ours) | $15.8 \pm 0.6$ | 26.2 |
| Contrastive (Ours) | $16.1 \pm 0.6$ | 26.9 |
| Masked + Contrastive (Ours) | $\mathbf{16.4 \pm 0.6}$ | **26.0** |

Table 3. Ablation study of masked contrastive pre-training on scene graph to image generation.

For qualitative comparison, we provide image samples with different scene graph embeddings in Fig. 5. Compared with one-hot embeddings and scene layout representations, we observe that embeddings from masked pre-training and contrastive pre-training enable SGDiff to generate images that are more realistic and comply better with the scene graphs. Compared to contrastive pretraining, we observe that embeddings obtained from masked pretraining tend to produce images that capture local structures better with
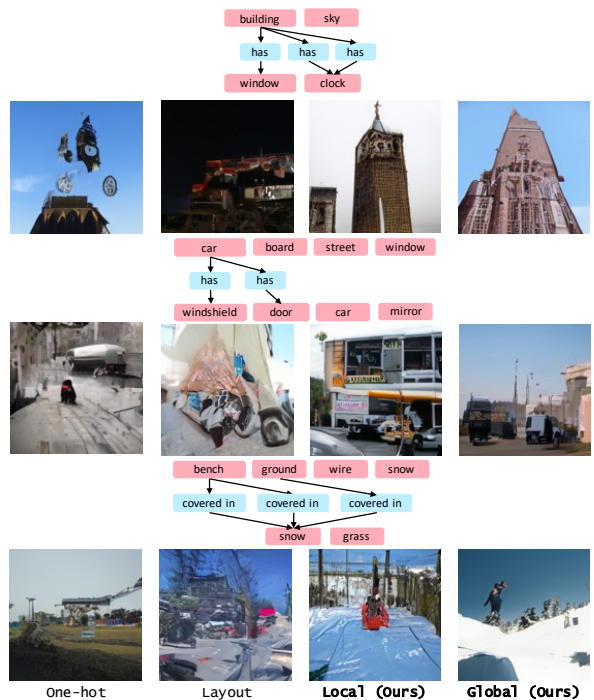


Figure 5. Latent diffusion with different SG embeddings.

8

more fine-grained object details, whereas embeddings from contrastive pretraining tend to focus more on matching the global structures at the overall image level.

## 5. Conclusion

This paper proposes a new framework, SGDiff, for image generation from scene graphs. SGDiff uses a masked contrastive pre-training approach to obtain scene graph emebeddings that allow for improved alignment between scene graphs and images, while also leveraging latent diffusion for improved scalability and generation quality. As a result, SGDiff produces more realistic and compliant images than previous methods that rely on manually crafted scene graph representations, such as scene layouts. SGDiff also makes image generation semantically controllable, allowing for easier manipulation of images through scene graph editing. Evaluation on standard datasets such as Visual Genome and COCO-Stuff show that SGDiff outperforms state-of-the-art approaches both qualitatively and quantitatively.

## References

[1] Oron Ashual and Lior Wolf. Specifying object attributes and relations in interactive scene generation. In *IEEE International Conference on Computer Vision*, pages 4561–4569, 2019. 2, 5, 6, 7

[2] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 18208–18218, 2022. 4

[3] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. In *International Conference on Learning Representations*, 2021. 4

[4] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *International Conference on Learning Representations*, 2018. 1

[5] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Cocostuff: Thing and stuff classes in context. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1209–1218, 2018. 2, 5

[6] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022. 4

[7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*, pages 1597–1607, 2020. 4

[8] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. *arXiv preprint arXiv:2210.11427*, 2022. 1

[9] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. Vqgan-clip: Open domain image generation and editing with natural language guidance. *arXiv preprint arXiv:2204.08583*, 2022. 1

[10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019. 4

[11] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems*, volume 34, pages 8780–8794, 2021. 1, 2

[12] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. *arXiv preprint arXiv:2203.13131*, 2022. 1

[13] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 10696–10706, 2022. 2

[14] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 4

[15] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020. 4

[16] Sen He, Wentong Liao, Michael Ying Yang, Yongxin Yang, Yi-Zhe Song, Bodo Rosenhahn, and Tao Xiang. Context-aware layout to image generation with enhanced object appearance. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 15049–15058, 2021. 2

[17] Roei Herzig, Amir Bar, Huijuan Xu, Gal Chechik, Trevor Darrell, and Amir Globerson. Learning canonical representations for scene graph to image generation. In *European Conference on Computer Vision*, pages 210–227. Springer, 2020. 2, 5, 6, 7

[18] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, volume 30, 2017. 5

[19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851, 2020. 1, 2, 5

[20] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 1, 2

[21] Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. *J. Mach. Learn. Res.*, 6:695–709, 2005. 5

[22] Justin Johnson, Agrim Gupta, and Li Fei-Fei. Image generation from scene graphs. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1219–1228, 2018. 1, 2, 5, 6, 7, 12, 13, 14

[23] Justin Johnson, Ranjay Krishna, Michael Stark, Li-Jia Li, David Shamma, Michael Bernstein, and Li Fei-Fei. Image retrieval using scene graphs. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3668–3678, 2015. 1, 2

[24] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. *arXiv preprint arXiv:2210.09276*, 2022. 4

[25] Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2426–2435, 2022. 1, 2

[26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015. 5

[27] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 5

[28] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017. 3

[29] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*, 123(1):32–73, 2017. 1, 2, 5

[30] Yikang Li, Tao Ma, Yeqi Bai, Nan Duan, Sining Wei, and Xiaogang Wang. Pastegan: A semi-parametric method to generate image from scene graph. *Advances in Neural Information Processing Systems*, 32, 2019. 2, 5, 6, 7, 12, 13, 14

[31] Zejian Li, Jingyu Wu, Immanuel Koh, Yongchuan Tang, and Lingyun Sun. Image synthesis from layout with locality-aware mask adaption. In *IEEE International Conference on Computer Vision*, pages 13819–13828, 2021. 2

[32] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32, 2019. 4

[33] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations*, 2021. 1, 4

[34] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob Mcgrew, Ilya Sutskever, and Mark Chen. GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*, pages 16784–16804, 2022. 2

[35] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 4

[36] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *IEEE International Conference on Computer Vision*, pages 2085–2094, 2021. 4

[37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763, 2021. 2, 4

[38] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67, 2020. 2

[39] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 1, 2

[40] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in Neural Information Processing Systems*, 32, 2019. 1

[41] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, pages 1278–1286, 2014. 5

[42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 1, 2, 4

[43] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 5

[44] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *Special Interest Group on Computer Graphics and Interactive Techniques Conference Proceedings*, pages 1–10, 2022. 2

[45] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 1, 2, 4

[46] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, volume 29, 2016. 5

[47] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008. 3

[48] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265, 2015. 1, 2

[49] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Advances in*

*Neural Information Processing Systems*, volume 32, 2019. 1, 2, 5

[50] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2020. 1, 2, 5

[51] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vl-bert: Pre-training of generic visual-linguistic representations. In *International Conference on Learning Representations*, 2019. 4

[52] Wei Sun and Tianfu Wu. Image synthesis from reconfigurable layout and style. In *IEEE International Conference on Computer Vision*, pages 10531–10540, 2019. 2

[53] Tristan Sylvain, Pengchuan Zhang, Yoshua Bengio, R Devon Hjelm, and Shikhar Sharma. Object-centric image generation from layouts. In *The AAAI Conference on Artificial Intelligence*, pages 2647–2655, 2021. 2

[54] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2818–2826, 2016. 5

[55] Subarna Tripathi, Anahita Bhiwandiwalla, Alexei Bastidas, and Hanlin Tang. Using scene graph context to improve image generation. *arXiv preprint arXiv:1901.03762*, 2019. 2

[56] Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. In *Advances in Neural Information Processing Systems*, volume 34, pages 11287–11302, 2021. 2, 4

[57] Dani Valevski, Matan Kalman, Yossi Matias, and Yaniv Leviathan. Unitune: Text-driven image editing by fine tuning an image generation model on a single image. *arXiv preprint arXiv:2210.09477*, 2022. 4

[58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 2, 5

[59] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011. 5

[60] Tengfei Wang, Ting Zhang, Bo Zhang, Hao Ouyang, Dong Chen, Qifeng Chen, and Fang Wen. Pretraining is all you need for image-to-image translation. *arXiv preprint arXiv:2205.12952*, 2022. 1

[61] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmim: A simple framework for masked image modeling. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 9653–9663, 2022. 4

[62] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Yingxia Shao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *arXiv preprint arXiv:2209.00796*, 2022. 1, 2

[63] Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao.

Vinvl: Revisiting visual representations in vision-language models. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5579–5588, 2021. 4

[64] Bo Zhao, Weidong Yin, Lili Meng, and Leonid Sigal. Layout2image: Image generation from layout. *International Journal of Computer Vision*, 128(10):2418–2435, 2020. 2

# Supplementary Material

## A. Implementation Details

As illustrated in main text, our proposed framework SGDiff consists of three (pre-)training stages, masked contrastive pre-training for SG encoder, variational autoencoder pre-training for latent embedding of images, and latent diffusion training for SG-based image generation. Here we introduce the concrete network and optimization details of these stages in Tab. 4 and Tab. 5. We use the same settings on both VG and COCO-Stuff datasets.

In masked contrastive pre-training, the nodes and edges of SGs are all preprocessed into 512-dimensional vectors for SG encoder. The ViT model tokenizes input images with patch size of 32×32. Specifically, we set the ratio of random mask to 0.3 in masked autoencoding branch, and set the ratio between masked autoencoding loss and contrastive loss to 10:1 for facilitating the optimization. In variational autoencoder pre-training, we embed input images into compact latents with a downsampling factor of 8, and maximize the decoding ability by optimizing the MSE objective. And the ratio between KL divergence and MSE is set to 8:10. In latent diffusion training, we use cross-attention mechanism for conditional diffusion process in all experiments.

| Stage | Module | Backbone | Input | Output | Blocks |
|---|---|---|---|---|---|
| Masked Contrastive Pre-Training | Image Encoder | ViT-B/32 | 256×256×3 | 512 | 11 |
| | SG Encoder | Graph Convolution | 512 | 512 | 5 |
| Variational Autoencoder Pre-Training | Encoder | ResNet | 256×256×3 | 32×32×4 | 3 |
| | Decoder | ResNet | 32×32×4 | 256×256×3 | 3 |
| Latent Diffusion Training | UNet-Encoder | ResNet + CrossAttention | 32×32×4 | 8×8×1024 | 9 |
| | UNet-MiddleBlock | ResNet | 8×8×1024 | 8×8×1024 | 3 |
| | UNet-Decoder | ResNet + CrossAttention | 8×8×1024 | 32×32×4 | 10 |

Table 4. Network Details of SGDiff.

| Stage | Optimizer | Batch Size | LR | Loss |
|---|---|---|---|---|
| Masked Contrastive Pre-Training | Adam | 16 | $5.0\,e^{-4}$ | Masked Autoencoding : Contrastive = 10 : 1 |
| Variational Autoencoder Pre-Training | Adam | 128 | $1.0\,e^{-5}$ | KL Divergence : MSE = 8 : 10 |
| Latent Diffusion Training | Adam | 16 | $1.0\,e^{-6}$ | $L_2$ Distance |

Table 5. Optimization Details of SGDiff.

## B. More Synthesis Results

In previous qualitative evaluations, we have shown the synthesis comparison results on VG dataset. Here, we provide more results on COCO-Stuff dataset with the resolution of 256×256 in Fig. 6 and Fig. 7. SGDiff exhibits the superiority of generation quality on COCO-Stuff dataset, and SGDiff can generate images that are more realistic and semantic-compliant than previous methods Sg2Im [22] and PasteGAN [30]. The results not only reveal the SG embeddings learned by our masked contrastive pre-training are effective in graph-image semantic alignment, but also demonstrate the efficacy of our local-global conditional latent diffusion.
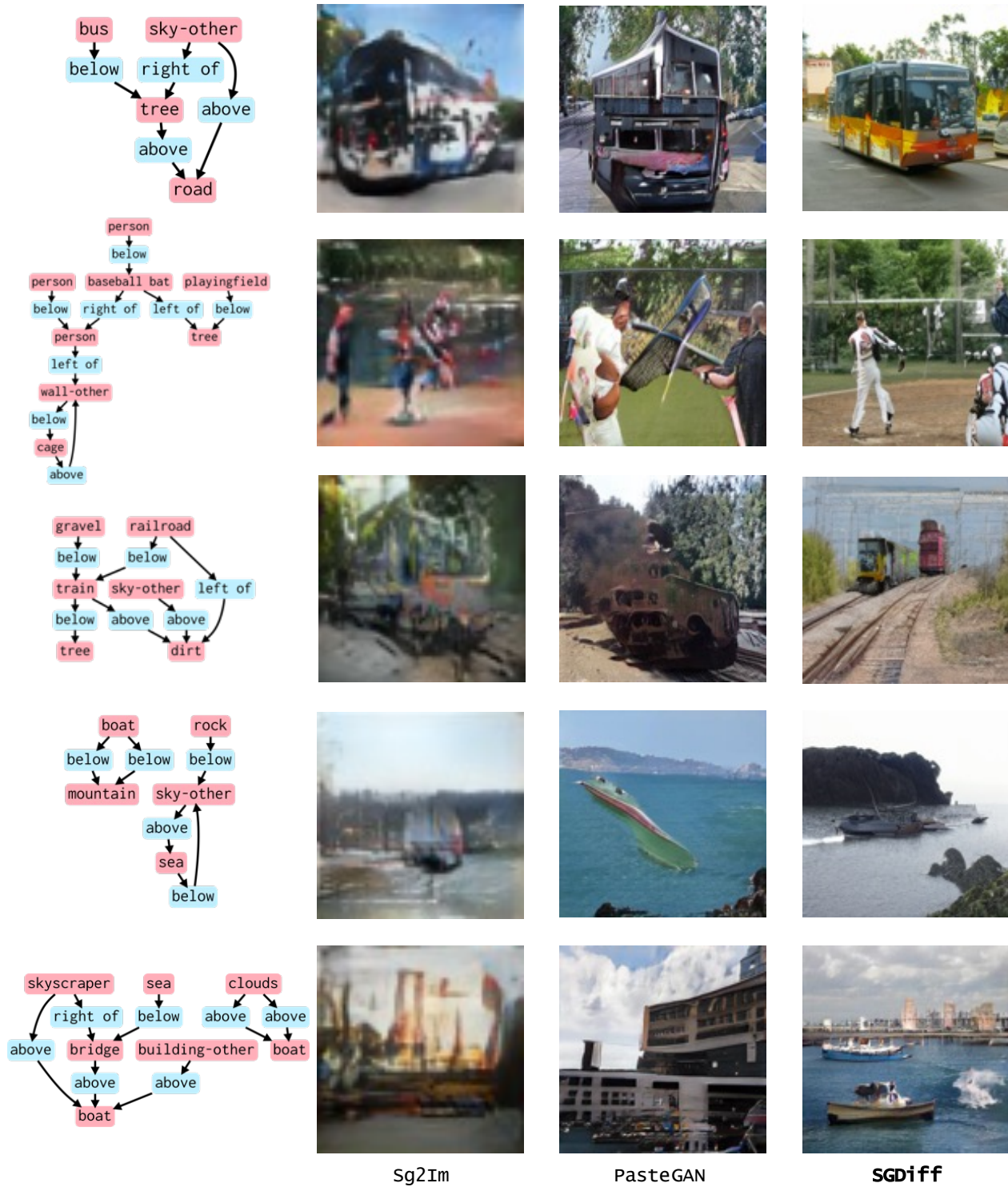
Figure 6. Image samples (256×256) generated by Sg2Im [22], PasteGAN [30], and our SGDiff given the same scene graphs.
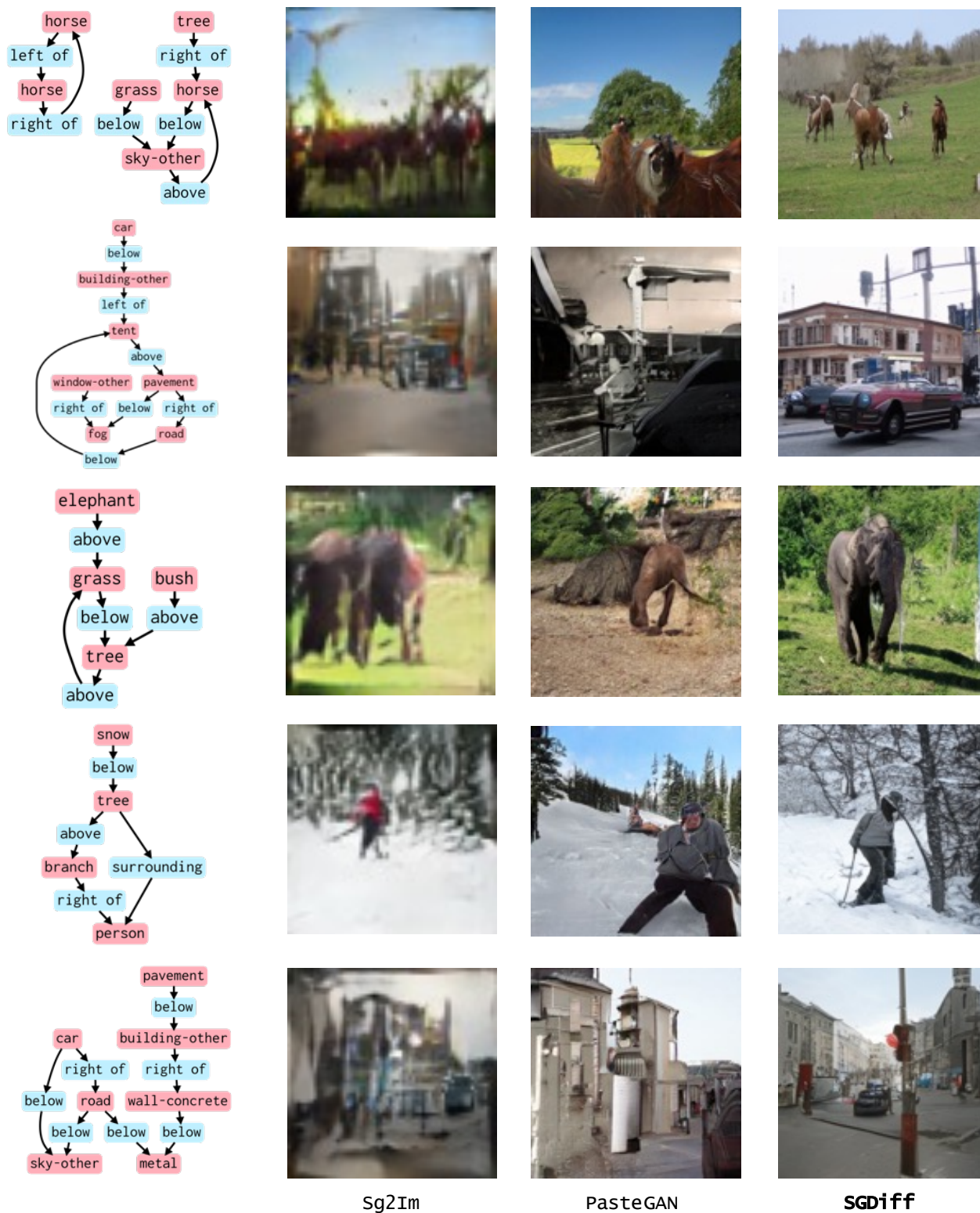
Figure 7. Image samples (256×256) generated by Sg2Im [22], PasteGAN [30], and our SGDiff given the same scene graphs.